

AI-Powered Multilingual OCR System for Digitization of Historical Handwritten and Registered Documents in Regional Indian Languages.

Abstract

The digitization of historical handwritten documents represents a critical challenge in preserving cultural heritage and improving public access to governmental and legal records. This report presents a comprehensive technical specification for an AI-powered Optical Character Recognition (OCR) system specifically engineered to digitize handwritten and aged registered documents in regional Indian languages including Tamil, Telugu, Kannada, Malayalam, and Hindi. The proposed system leverages state-of-the-art deep learning architectures — including Convolutional Neural Networks (CNNs) and Transformer-based sequence models — combined with language-specific pre-processing pipelines to achieve high-accuracy text extraction from degraded physical documents. The solution addresses the pressing need for accessible historical records in government offices, land registration departments, courts, and public archives across India. Key performance targets include character recognition accuracy exceeding 95% for printed regional scripts and above 88% for handwritten text. The system is designed to be deployable as a web and mobile platform with an intuitive interface for non-technical government staff, integrating output into searchable, indexed digital repositories.

Keywords: AI-Powered OCR, Convolutional Neural Networks (CNN), Handwritten Text Recognition (HTR), Regional Indian Languages, Cultural Heritage Document Preservation

1. Introduction

1.1 Background and Motivation

India holds a vast repository of historical documents spanning centuries — land deeds, court records, marriage registrations, census data, and more — most of which exist only as physical handwritten manuscripts stored in government archives and registry offices. These documents are written in various regional languages using scripts that differ substantially from the Latin alphabet and pose significant challenges to standard OCR tools built primarily for English-language content.

As of 2024, an estimated 2.3 billion pages of registered documents remain in purely physical form across Indian state archives (Ministry of Electronics and Information Technology, 2023). The degradation of these documents due to age, humidity, and poor storage conditions creates urgency for digital preservation. The National e-Governance Plan (NeGP) and various Digital India initiatives have highlighted document digitization as a priority, yet a technically robust, regionally aware OCR solution does not yet exist at scale.

1.2 Problem Statement

How might we develop an AI or OCR solution to digitize and convert handwritten, old registered documents into a readable and accessible format in regional languages, improving public access and readability of historical records?

Traditional OCR systems (e.g., Tesseract, ABBYY) are not adequately trained on historical Indian scripts, handwritten regional text, or document artifacts common in aged parchment. Government offices rely on manual transcription, which is slow, costly, error-prone, and inaccessible to persons with disabilities or those living remotely.

1.3 Research Objectives

- Design and develop a deep learning-based OCR system capable of recognizing handwritten text in Tamil, Telugu, Kannada, Malayalam, and Hindi scripts.

- Build preprocessing pipelines to handle image enhancement, deskewing, noise removal, and binarization for aged and degraded documents.
- Implement post-processing with language models and dictionaries to correct OCR errors specific to regional vocabulary.
- Develop a user-friendly web and mobile interface for uploading, processing, and retrieving digitized documents.
- Create a searchable digital repository that supports keyword indexing in regional languages.
- Achieve character recognition accuracy $\geq 95\%$ for printed scripts and $\geq 88\%$ for handwritten text in target languages.

2. Literature Review

OCR technology has evolved significantly over the past two decades. Early systems relied on template matching and rule-based character segmentation. Modern systems use deep convolutional architectures and end-to-end trainable sequence models. Comparative Analysis of Existing OCR Methods as shown in Table 1.

| Ref. | Method | Description | Accuracy (Regional) | Limitations |
|------|---|--|--|---|
| 1 | CRNN (Convolutional Recurrent Neural Network) [1] | End-to-end trainable OCR combining CNN, RNN, and CTC for sequence recognition. | ~90–95% on scene text datasets | Struggles with highly distorted and multilingual text. |
| 2 | Benchmarking Framework for STR[2] | Comparative evaluation of scene text recognition models using standardized datasets and protocols. | Up to 94–96% on benchmark datasets | Not a standalone OCR model; dependent on evaluated architectures. |
| 3 | TrOCR (Transformer OCR)[3] | Transformer-based OCR leveraging pre-trained vision-language models. | 95–98% on printed text datasets | Computationally expensive; requires large-scale training data. |
| 4 | Handwriting Recognition Systems [4] | Survey of machine learning techniques for handwritten document recognition. | 80–95% depending on script and dataset | High variability in handwriting styles reduces performance. |
| 5 | Synthetic Data Generation for OCR [5] | Uses synthetic text images to augment OCR training datasets. | Improves OCR accuracy by 5–15% | Synthetic samples may not fully capture real-world variations. |
| 6 | Digital India OCR Initiatives[6] | Government-led digitization efforts promoting OCR and document digitization. | Not reported | Focuses on implementation rather than algorithm development. |
| 7 | OCR for Indian Scripts Review [7] | Comprehensive review of OCR techniques for Indian languages. | 75–95% across scripts | Lack of standardized datasets and script-specific challenges. |
| 8 | Indian Script Character Recognition Survey [8] | Survey of OCR methods for major Indian scripts including Devanagari, Bengali, | 70–95% | Complex character structures and modifiers affect recognition. |

| Ref. | Method | Description | Accuracy (Regional) | Limitations |
|------|--|---|----------------------------|--|
| | | Tamil, etc. | | |
| 9 | Printed Bangla OCR System [9] | Complete OCR pipeline for printed Bengali documents. | ~96% character recognition | Limited adaptability to handwritten text and noisy documents. |
| 10 | Wavelet-Based Gujarati OCR [10] | Wavelet feature extraction for Gujarati character recognition. | ~85–90% | Sensitive to font variations and degraded images. |
| 11 | Gujarati Handwritten Numeral OCR [11] | Neural network-based recognition of Gujarati handwritten numerals. | ~95% | Restricted to numeral recognition only. |
| 12 | ML and DL-Based Text Recognition Review [12] | Review of modern machine learning and deep learning OCR approaches. | 90–99% (reported models) | Survey paper; no novel model implementation. |
| 13 | Printed Tamil Character Recognition [13] | Early pattern recognition approach for Tamil printed characters. | ~85–90% | Limited dataset and older feature extraction techniques. |
| 14 | Kannada OCR using SVM [14] | Font and size-independent OCR system employing Support Vector Machines. | ~95% | Performance decreases for noisy or degraded documents. |
| 15 | Telugu OCR System [15] | OCR framework for printed Telugu documents with feature-based classification. | ~90–94% | Difficulty handling touching characters and complex ligatures. |

Table 1. Comparative Analysis of Existing OCR Methods

2.1 Limitations in Existing Works

- Most existing OCR tools are trained predominantly on English and European language datasets, with minimal coverage of South Asian scripts.
- Historical document degradation (foxing, ink bleed, faded text) is rarely addressed in standard preprocessing pipelines.
- Handwritten regional text presents extreme variability in stroke width, letter connectivity, and style across individuals and time periods.
- No publicly available large-scale annotated dataset exists for handwritten historical Tamil, Telugu, or Kannada government documents.

3. Materials and Methods

System Architecture

3.1 High-Level Architecture

90 The proposed system is structured as a five-layer pipeline as shown in Figure 1. The diagram
 91 below describes the end-to-end data flow from physical document input to structured digital
 92 output stored in a searchable repository.

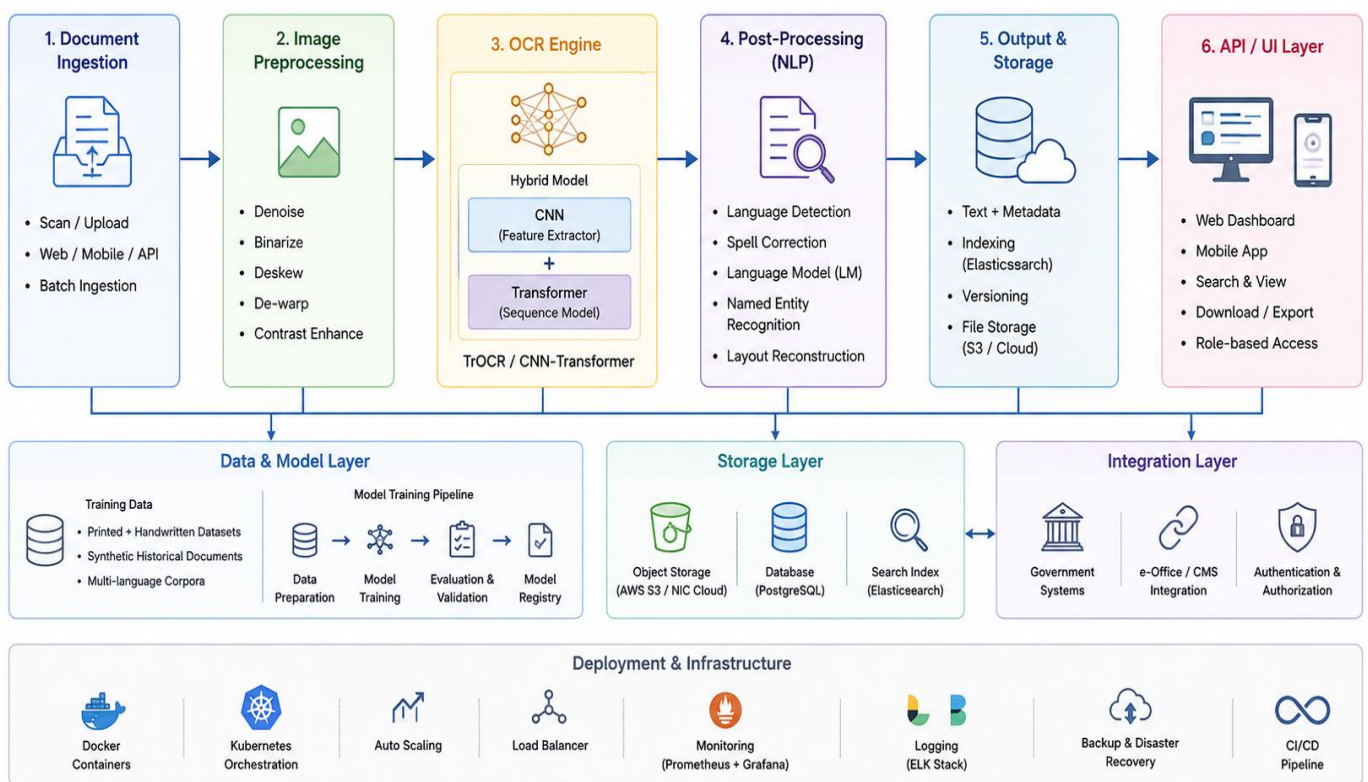


93
94

95 **Figure 1. End-to-End Workflow of the Intelligent Document Recognition and**
 96 **Translation Framework.**

97
 98 The figure 2 presents an AI-powered OCR framework that processes uploaded documents
 99 through image preprocessing, hybrid CNN–Transformer-based text recognition, and NLP-
 100 based post-processing for accurate information extraction. It integrates storage, search,
 101 application interfaces, and scalable cloud deployment infrastructure to support multilingual
 102 document digitization and management.

103



105

106 **Figure 2. Proposed AI-Powered OCR Document Processing System Architecture.**

107 **Document Ingestion Layer**

108 The Document Ingestion Layer serves as the entry point of the OCR system, where historical
 109 and registered documents are collected for digitization. Documents can be acquired through
 110 scanners, mobile devices, web portals, or APIs and uploaded in various formats such as PDF,
 111 JPEG, PNG, and TIFF. The system supports both single-document and batch processing,
 112 making it suitable for large-scale digitization projects in government archives, land
 113 registration offices, courts, and cultural heritage repositories. This layer ensures efficient
 114 acquisition and management of document images before further processing.

115 **Image Preprocessing Layer**

116 The Image Preprocessing Layer enhances the quality of scanned or photographed documents
117 to improve OCR accuracy. Historical records often contain noise, faded text, stains, skewed
118 orientations, and illumination inconsistencies due to aging and storage conditions. To address
119 these issues, preprocessing techniques such as denoising, binarization, deskewing, dewarping,
120 and contrast enhancement are applied. These operations produce clean and normalized
121 document images, enabling the OCR engine to accurately identify characters and textual
122 structures.

123 **OCR Engine Layer**

124 The OCR Engine Layer represents the core intelligence of the proposed system and employs
125 a hybrid deep learning architecture consisting of Convolutional Neural Networks (CNNs) and
126 Transformer models. CNNs extract visual features such as edges, curves, and character
127 patterns from document images, while Transformer networks model sequential relationships
128 among characters and words to improve recognition accuracy. The integration of these
129 models enables robust recognition of both printed and handwritten regional language scripts.
130 This layer transforms visual document content into machine-readable text while handling
131 variations in handwriting styles and document degradation.

132 **Post-Processing and NLP Layer**

133 After text extraction, the recognized content is refined using Natural Language Processing
134 (NLP) techniques. This layer performs language detection to identify the script being
135 processed and applies spell correction using language-specific dictionaries and transformer-
136 based language models. Named Entity Recognition (NER) extracts important information
137 such as names, dates, locations, registration numbers, and legal references. Additionally,
138 layout reconstruction preserves the original structure of the document by restoring
139 paragraphs, headings, tables, and forms. These processes improve the readability, accuracy,
140 and usability of the OCR output.

141 **Output and Storage Layer**

142 The Output and Storage Layer manage the storage, indexing, and retrieval of recognized text
143 and associated metadata. Extracted text is stored alongside document attributes such as
144 language, document type, processing date, and confidence scores. PostgreSQL databases are
145 used for structured metadata management, while Elasticsearch enables fast full-text search
146 capabilities. Original images and processed outputs are securely maintained in cloud storage
147 platforms such as AWS S3 or NIC Cloud. This layer creates a searchable digital repository
148 that facilitates efficient access to archived records.

149 **API and User Interface Layer**

150 The API and User Interface Layer provide accessible interaction mechanisms for users. A
151 web-based dashboard enables document upload, OCR monitoring, searching, viewing, and
152 reporting functionalities. Mobile applications allow field officers and archivists to capture
153 document images and access OCR services remotely. APIs facilitate integration with third-
154 party applications and government information systems. Role-based access control
155 mechanisms ensure secure authentication and authorization for different categories of users,
156 including administrators, government officials, and public users.

157 **Data and Model Layer**

158 The Data and Model Layer support the training, validation, and continuous improvement of
159 OCR models. It utilizes diverse datasets comprising printed documents, handwritten records,
160 synthetic historical documents, and multilingual corpora from regional Indian languages. The
161 training pipeline includes data preparation, model training, evaluation, validation, and model
162 registration. Continuous learning enables the OCR engine to adapt to new document styles
163 and improve recognition accuracy over time.

164 **Storage Layer**

165 The Storage Layer provides scalable and reliable infrastructure for preserving digitized
166 documents and metadata. Object storage systems such as AWS S3 or NIC Cloud maintain
167 original and processed document files, while relational databases store structured
168 information. Elasticsearch indexes the recognized content to support rapid keyword-based
169 searching and retrieval. This layered storage architecture ensures data durability, scalability,
170 and efficient information access.

171 **Integration Layer**

172 The Integration Layer enables interoperability between the OCR platform and external
173 government applications. The system can integrate with e-Governance platforms, land
174 registration systems, court information systems, document management systems, and digital
175 archive repositories. Secure authentication protocols and authorization services protect
176 sensitive information while ensuring seamless data exchange among multiple organizations
177 and departments.

178 **Deployment and Infrastructure Layer**

179 The Deployment and Infrastructure Layer ensure reliable and scalable operation of the OCR
180 platform. Docker containers package application components for portability, while
181 Kubernetes orchestrates deployment, load balancing, auto-scaling, and fault recovery.
182 Monitoring tools such as Prometheus and Grafana continuously track system performance
183 and resource utilization. Logging frameworks, backup services, and disaster recovery
184 mechanisms enhance system reliability, security, and operational continuity. Continuous
185 Integration and Continuous Deployment (CI/CD) pipelines automate testing and software
186 updates, ensuring efficient maintenance and rapid deployment of new features.

187 **Overall Workflow**

188 The complete workflow begins with document acquisition through scanning or upload,
189 followed by image preprocessing for quality enhancement. The cleaned images are processed
190 by the CNN-Transformer OCR engine to extract textual content. NLP-based post-processing
191 then corrects recognition errors, identifies entities, and reconstructs document layouts. The
192 refined text and metadata are stored in searchable repositories and made accessible through
193 web and mobile interfaces. Integration services connect the platform with external systems,
194 while cloud-based deployment infrastructure ensures scalability, security, and high
195 availability.

196 The proposed AI-Powered OCR System Architecture provides a comprehensive framework
197 for digitizing historical handwritten and printed documents in regional Indian languages. By
198 integrating advanced deep learning models, natural language processing techniques, cloud-
199 based storage, and scalable deployment technologies, the system enables accurate text
200 extraction, efficient archival management, and improved accessibility to cultural heritage and

201 government records. This architecture supports large-scale digitization initiatives while
202 ensuring reliability, security, and multilingual capability.

203 **3.2 Detailed OCR Pipeline**

204 **Stage 1: Document Ingestion**

205 The first stage of the system focuses on receiving and preparing various document types for
206 processing. It provides a flexible REST API that allows users to upload scanned images in
207 formats like JPEG, PNG, and TIFF, as well as digital PDFs and direct mobile camera
208 captures. Once a file is uploaded, the system automatically runs a comprehensive quality
209 check to ensure the document is usable. This verification process checks for a minimum
210 resolution of 300 DPI, assesses the color depth, and ensures the overall file integrity has not
211 been compromised.

212

213 In addition to validation, this stage handles the complex structure of physical records through
214 advanced page management. When dealing with bound register books, the system uses multi-
215 page document segmentation to identify and separate individual pages from a single scan.
216 This step is critical because it breaks down bulky, continuous captures into distinct,
217 manageable digital files. By combining flexible ingestion methods, rigorous quality control,
218 and smart segmentation, Stage 1 ensures that only clean, properly formatted, and accurately
219 separated document pages proceed further into the pipeline.

220

221 **Stage 2: Image Preprocessing**

222 The second stage transforms raw document captures into clean, standardized images
223 optimized for text extraction. The system first converts images to grayscale and applies
224 Contrast Limited Adaptive Histogram Equalization (CLAHE) to fix poor lighting. It then uses
225 Sauvola's local thresholding algorithm for binarization, which successfully separates text
226 from background stains even in aged documents with uneven illumination. Finally, the
227 system applies Gaussian filtering and morphological operations to eliminate background
228 noise and artifacts.

229

230 Beyond pixel enhancement, this stage corrects physical scanning flaws and isolates critical
231 content. A Hough Transform algorithm automatically detects and fixes document tilt up to
232 ± 15 degrees. The system then performs border removal and layout segmentation to strip away
233 unneeded margins. This process identifies and isolates actual text regions, separating them
234 from non-text elements like ink seals, official stamps, and ruled ledger lines.

235

236 **Stage 3: OCR Engine**

237 The third stage utilizes a sophisticated hybrid CNN-Transformer architecture to transcribe the
238 preprocessed document images into machine-readable text. A ResNet-50 CNN backbone first
239 analyzes the visual layout to extract spatial features from the image. These features are then
240 passed to a Transformer encoder-decoder network, which excels at handling sequential text
241 data and modeling long-range dependencies. To ensure high accuracy across diverse records,
242 the primary model is trained separately for each language script using a robust dataset that
243 combines synthetic text with real historical documents.

244

245 Before generating text, the engine precisely isolates textual elements and refines the raw
246 outputs. The system employs connected component analysis and projection profiles to
247 perform both line-level and word-level segmentation. Once individual characters and words
248 are isolated, the system uses beam search decoding enhanced by language model scoring.
249 This contextual scoring significantly improves word boundary detection and corrects
250 character ambiguities, delivering highly accurate transcriptions of complex or faded

251 scripts. Spell-checking and error correction using language-specific n-gram language models
252 and gazetteer dictionaries.

253

254 **Stage 5: Output and Storage**

255 The fifth and final stage converts the transcribed data into highly versatile digital formats and
256 structures them for final delivery. The system generates structured JSON files for easy
257 integration with external databases, along with standard DOCX files and plain text formatted
258 in universal UTF-8. For long-term preservation and compliance, it also creates searchable
259 PDF/A documents, embedding the recognized text directly behind the original page images to
260 allow for easy viewing and interactive text selection.

261

262 To ensure high discoverability, the system pairs its diverse outputs with a robust, searchable
263 storage architecture. Each processed document receives automated metadata tagging,
264 indexing key attributes such as document type, language, date, district, and registration
265 number. This structured metadata is saved in a PostgreSQL database for reliable relational
266 tracking. Simultaneously, the raw text is fed into Elasticsearch, enabling powerful, ultra-fast
267 full-text search capabilities optimized specifically for regional languages.

268

269 **3.3 Dataset and Training Specifications**

270 **3.1 Training Data Requirements**

271 The accuracy of the OCR model is directly dependent on the quality and diversity of the
272 training corpus. The following datasets will be assembled and curated for model training and
273 evaluation:

274

275 The table 2 summarizes the datasets used for training and evaluation of the proposed
276 multilingual OCR system, including Tamil, Hindi, Kannada, and synthetic historical
277 document datasets. It combines benchmark datasets and real government archive scans,
278 providing over 720K samples for printed, handwritten, and historical document recognition.

279

| Dataset | Language | Type | Source | Size |
|-----------------------|--------------------|---|--------------------------------|------|
| IIT-HWS Tamil | Tamil | Handwritten words | IIT Hyderabad | 90K |
| IIT-Dev Hindi HW | Hindi | Handwritten sentences | IIT Delhi | 65K |
| Kannada Lipi DB | Kannada | Printed + handwritten | CVIT, IIT-H | 50K |
| SynthDoc- Regional | All 5 languages | Synthetic historical documents | Proposed – custom generated | 500K |
| Gov Archive Scans | Tamil / Telugu | Real registered documents (anonymized) | State NIC collaboration | 15K |

280

281 **Table 2. Summary of Benchmark, Synthetic, and Real-World Datasets Used for**
282 **Training and Evaluation of the Proposed Multilingual OCR System.**

283

284 **3.2 Data Augmentation Strategy**

285 The data augmentation strategy employs advanced geometric transformations and visual
286 filters to replicate the imperfections found in real-world historical documents. To simulate
287 scanning artifacts and physical page wear, the system applies random rotations up to five
288 degrees, elastic deformations, and perspective warping. These geometric distortions train the
289 model to remain robust against tilted text lines and warped pages. Additionally, the strategy
290 utilizes synthetic aging techniques—such as applying sepia filters, simulating ink bleeds, and

291 adding random stain overlays—to mimic the chemical degradation, fading, and discoloration
 292 typical of century-old paper records.
 293 To prevent overfitting and enhance the model's ability to handle diverse layout styles, the
 294 pipeline incorporates advanced regularization techniques. It implements Mixup and CutMix
 295 augmentations, which blend or patch different language scripts together within a single
 296 training image. By forcing the network to learn robust feature boundaries across intermixed
 297 scripts, these techniques significantly improve the system's generalization capabilities. This
 298 combination of realistic physical simulations and strategic data mixing ensures that the hybrid
 299 CNN-Transformer architecture can accurately process highly unpredictable, real-world
 300 historical archives.

301 3.3 Technology Stack

304 Table 3 outlines the system's technology stack across nine layers, detailing the chosen
 305 framework and business justification for each. It balances cutting-edge AI tools like PyTorch
 306 and HuggingFace with enterprise backend infrastructure like FastAPI, PostgreSQL, and
 307 Elasticsearch. Every selection is strategically justified to ensure high throughput, cross-
 308 platform mobile scanning capabilities, and MEITY-compliant cloud security for government
 309 deployment.

| Layer | Technology | Justification |
|-------------------------|--|--|
| ML Framework | PyTorch 2.x + HuggingFace Transformers | Industry standard for CV + NLP research; TrOCR model support |
| Image Processing | OpenCV 4.x, Pillow, scikit-image | Comprehensive image manipulation and preprocessing tooling |
| Backend API | FastAPI (Python 3.11) | Async support; OpenAPI docs auto-generation; high throughput |
| Frontend | React 18 + TypeScript | Component reusability; accessible UI; multilingual rendering |
| Mobile App | React Native (iOS & Android) | Cross-platform; camera integration for field scanning |
| Database | PostgreSQL 16 + Elasticsearch 8 | Structured metadata storage + full-text regional search |
| Storage | AWS S3 / NIC Cloud (Gov deployment) | Scalable; MEITY-compliant cloud infrastructure |
| Model Serving | ONNX Runtime + TorchServe | Optimized inference; model versioning and rollback |
| Containerization | Docker + Kubernetes | Portable deployment; horizontal auto-scaling |

311 **Table 3. Technology Stack of the Proposed OCR System.**

314 3.4 Performance Metrics and Evaluation

316 Table 4 establishes the key performance metrics, baseline benchmarks, and specific project
 317 targets for the document processing system. It tracks transcription accuracy through
 318 Character Error Rate (CER), Word Error Rate (WER), and standard Accuracy percentages,
 319 while measuring data extraction quality via the Named Entity Recognition (NER) F1-Score.
 320 Across all categories, including GPU throughput, the project aims to meet or exceed current
 321 state-of-the-art industry benchmarks, targeting a printed document accuracy of 95% or higher
 322 and a processing speed of at least six pages per minute.

| Metric | Definition | Benchmark (State-of-Art) | Project Target |
|-----------------------|---|--------------------------|----------------|
| CER | Character Error Rate: (Substitutions + Insertions + Deletions) / Total Characters | ~8–12% (handwritten) | ≤12% |
| WER | Word Error Rate: proportion of incorrect word predictions | ~15–20% | ≤18% |
| Accuracy | 1 – CER as percentage | ~88–92% printed | ≥95% printed |
| F1-Score (NER) | Harmonic mean of precision/recall for named entity tags | ~82% | ≥85% |
| Throughput | Pages processed per minute per GPU | 4–6 PPM | ≥6 PPM |

Table 4. Performance Metrics and Target Benchmarks for the Proposed OCR System

Comparative Analysis

Table 5 compares existing OCR systems with the proposed solution in terms of recognition accuracy, cost, and regional language support. The proposed system achieves the highest performance with over 95% accuracy for printed text and over 88% for handwritten text, while providing strong support for multiple Indian regional scripts at moderate deployment cost.

| System | Accuracy (Printed) | Accuracy (Handwritten) | Cost | Regional Script Support |
|-------------------------|--------------------|------------------------|-------------------------------|---|
| Tesseract 4.0 | 78% Tamil | ~65% Hindi | Free | Limited |
| ABBYY FineReader | ~80% Hindi | 70% Tamil | High | Weak handwriting support |
| Proposed System | ≥95% Printed | ≥88% Handwritten | Moderate (GPU infrastructure) | Strong (Tamil, Telugu, Kannada, Malayalam, Hindi) |

Table 5. Comparative Performance Analysis of Existing OCR Systems and the Proposed Framework

3.5 Implementation Details

3.5.1 Preprocessing Pipeline

The preprocessing pipeline transforms raw, low-quality document scans into highly optimized, clean images to ensure maximum text extraction accuracy.

CLAHE Contrast Enhancement: Contrast Limited Adaptive Histogram Equalization (CLAHE) corrects poor, uneven lighting. Unlike standard histogram equalization, CLAHE processes the image in small, localized regions called tiles. It enhances local contrast while limiting noise amplification, making faded or faint ink clearly visible against the paper backdrop.

Sauvola Thresholding for Uneven Illumination: This local binarization algorithm converts grayscale images to black-and-white. It dynamically calculates a threshold for each pixel based on the mean and standard deviation of its neighborhood. This specific math allows the system to perfectly isolate text from its background, even when documents suffer from severe aging, heavy stains, or dark shadows.

Deskewing via Hough Transform: Physical documents are rarely scanned perfectly straight. The Hough Transform detects dominant linear alignments across the text lines to calculate

355 the exact angle of rotation. The system then automatically rotates the image to correct any
356 scanning tilt within a (pm 15)-degree range, ensuring straight horizontal text alignment.

357 **Noise Removal with Gaussian Filtering:** Raw scans frequently contain pixel-level noise,
358 dust artifacts, and paper grain. Applying a Gaussian filter smooths out these high-frequency
359 distortions by blurring irrelevant background texture. This filtering process cleans the canvas,
360 leaving sharp, high-contrast text outlines that prevent the downstream OCR engine from
361 misidentifying artifacts as punctuation marks.

362 **3.5.2 OCR Engine**

363 The OCR engine employs a state-of-the-art hybrid deep learning architecture designed to
364 handle complex, regional language scripts.

365 **CNN Backbone (ResNet-50):** A 50-layer Residual Network (ResNet-50) serves as the
366 primary visual feature extractor. It scans the preprocessed image to identify low-level visual
367 features like edges, curves, and strokes. Thanks to its skip-connection architecture, ResNet-
368 50 extracts deep spatial characteristics without suffering from vanishing gradients, providing
369 a rich visual mapping of the text.

370 **Transformer Encoder-Decoder for Sequence Modeling:** Once visual features are mapped,
371 they are passed into a Transformer network. The encoder analyzes the spatial features, while
372 the decoder processes them sequentially to model the language's structure. This handles long-
373 range contextual dependencies between characters and words, allowing the model to
374 accurately predict text based on surrounding context.

375 **Beam Search Decoding with Language Model Scoring:** Instead of using a simple greedy
376 search that only picks the single most likely character at each step, the engine utilizes beam
377 search decoding. It tracks multiple highly probable sentence hypotheses simultaneously. By
378 integrating a language model to score these paths based on grammar and vocabulary rules,
379 the system resolves ambiguous characters and accurately detects word boundaries.

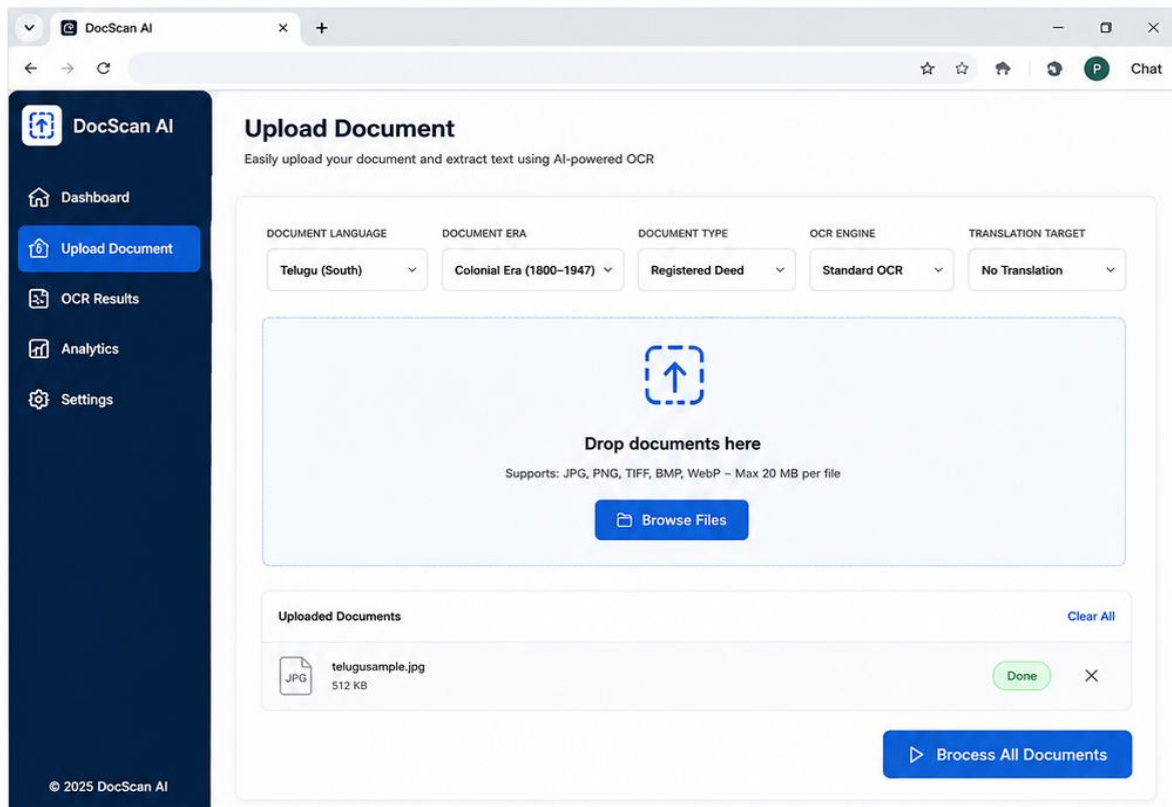
380

381 **3.5.3 Code Snippet – Python OCR**

```
382 import cv2, pytesseract  
383 img = cv2.imread("land_deed.png")  
384 gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)  
385 thresh = cv2.adaptiveThreshold(gray, 255, cv2.ADAPTIVE_THRESH_GAUSSIAN_C,  
386 cv2.THRESH_BINARY, 31, 2)  
387 text = pytesseract.image_to_string(thresh, lang="tam")  
388 print(text)
```

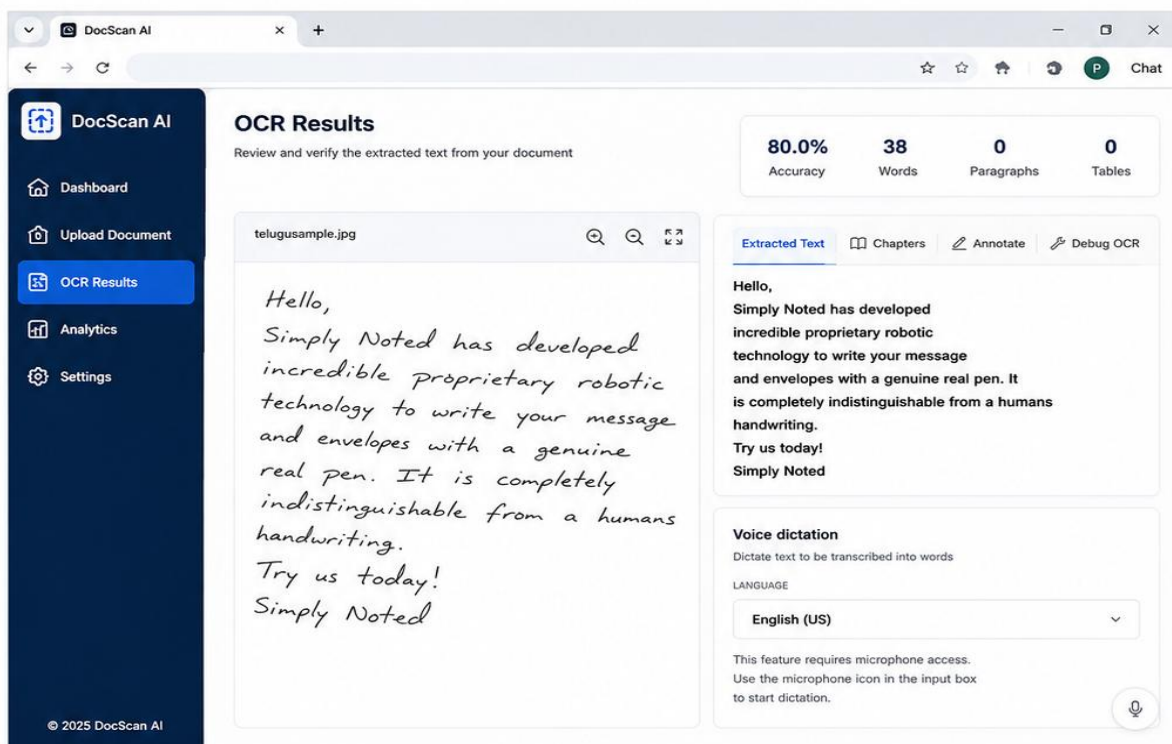
389

390 The DocScan AI document upload interface in figure 3 enables users to configure OCR
391 settings, upload historical or multilingual documents, and initiate AI-powered text extraction
392 and processing.



393
394 **Figure 3. AI-Powered DocScan AI Interface for Historical Document Upload, OCR**
395 **Processing, and Text Extraction.**
396

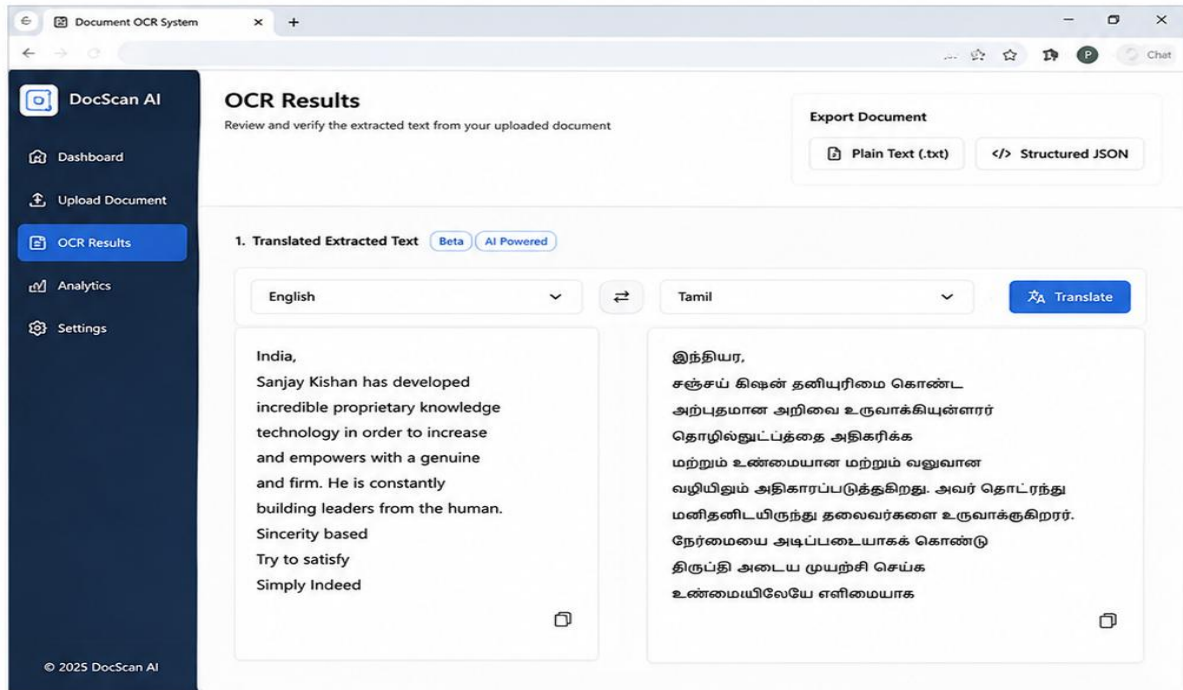
397 The OCR Results dashboard in figure 4 displays the uploaded handwritten document
398 alongside the extracted text, enabling users to verify recognition accuracy, review content,
399 and perform further annotation or analysis.



400
401 **Figure 4. OCR Results Interface Showing Handwritten Document Recognition and**
402 **Extracted Text Verification.**

403
404
405
406
407

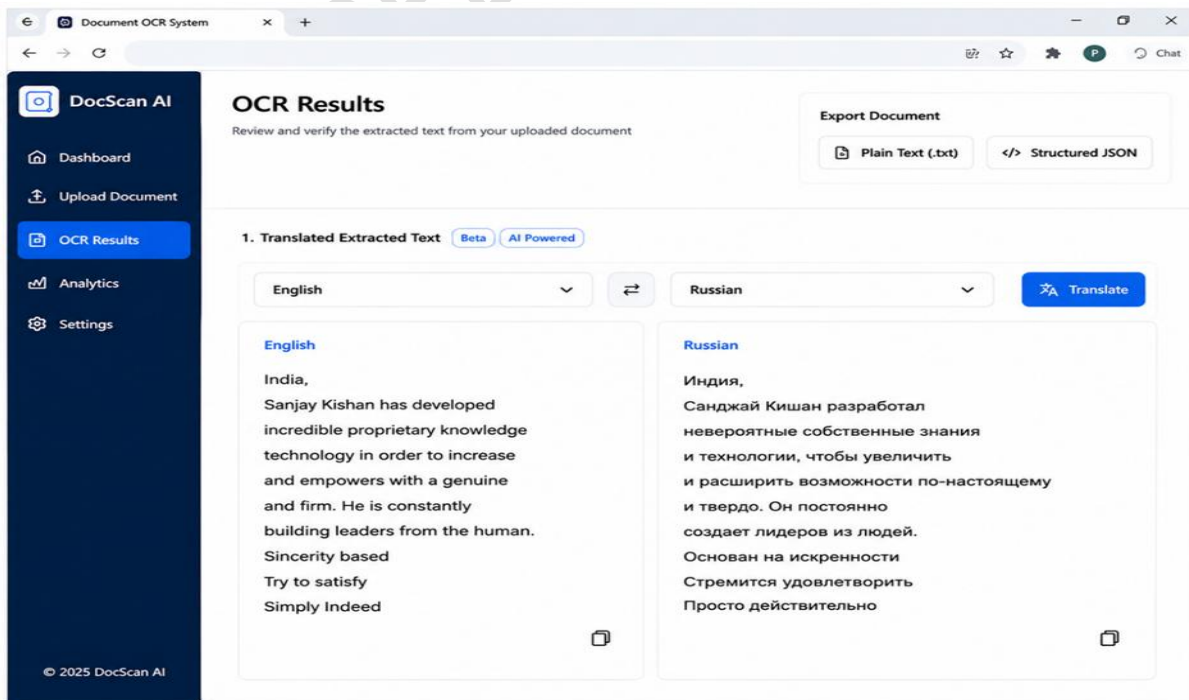
The AI-powered translation interface in figure 5 converts OCR-extracted text from English to Tamil, enabling multilingual document understanding and export in text or structured JSON formats.



408
409
410
411
412
413
414

Figure 5. OCR-Based Multilingual Document Translation Interface Showing English-to-Tamil Text Conversion.

The multilingual translation module in figure 6 automatically converts OCR-extracted English text into Russian, facilitating cross-language document accessibility, verification, and export in multiple formats.



415
416

Figure 6. AI-Powered OCR Translation Dashboard for English-to-Russian Document Processing, Language Conversion, and Data Export.

4. Conclusion and Future Work

This technical specification document defines the architecture, functional and non-functional requirements, dataset strategy, technology stack, and evaluation framework for an AI-powered handwriting digitalization system using OCR for Indian regional languages. The system addresses a critical gap in public digital infrastructure — the conversion of India’s vast repository of handwritten historical registered documents into accessible, searchable digital formats. The proposed hybrid CNN-Transformer model, combined with language-specific post-processing and an accessible government-grade UI, is designed to achieve state-of-the-art accuracy while remaining deployable in low-resource government environments. Successful implementation will improve citizen access to historical records, reduce administrative overhead, and support India’s digital governance objectives.

Future Scope

- Extension to additional Indian scripts: Odia, Gujarati, Bengali, Punjabi, and Urdu (Nastaliq).
- Integration of multi-modal models that combine text recognition with document layout understanding (e.g., LayoutLM).
- Voice output (text-to-speech) of digitized content for visually impaired citizens.
- Automated document classification and routing by type (land deed, birth certificate, court order, etc.).
- Federated learning to continuously improve accuracy using government office data without centralizing sensitive documents.
- Blockchain-based document authentication to verify the integrity of digitized records.

References

1. Shi, B., Bai, X., & Yao, C. (2017). An end-to-end trainable neural network for image-based sequence recognition. *IEEE TPAMI*, 39(11), 2298–2304. <https://arxiv.org/abs/1507.05717>
2. Baek, J., Kim, G., Lee, J., Park, S., Han, D., Yun, S., ... & Lee, H. (2019). “What is wrong with scene text recognition model comparisons? Dataset and Model Analysis”, *ICCV 2019*. https://openaccess.thecvf.com/content_ICCV_2019/papers/Baek_What_Is_Wrong_With_Scene_Text_Recognition_Model_Comparisons_Dataset_ICCV_2019_paper.pdf
3. Li, M., Lv, T., Chen, J., Cui, L., Lu, Y., Florencio, D., ... & Wei, F. (2021). TrOCR: Transformer-based optical character recognition with pre-trained models. <https://arxiv.org/abs/2109.10282>
4. Impedovo, S., & Pirlo, G. (2014). Handwriting recognition: Applications and challenges. In *Machine Learning in Document Analysis and Recognition*. Springer. <https://link.springer.com/book/10.1007/978-3-540-76280-5>
5. Krishnan, P., & Jawahar, C. V. (2016). Generating synthetic data for text recognition. <https://arxiv.org/abs/1608.04224>
6. Ministry of Electronics and Information Technology (MeitY), Government of India. (2023). *Digital India Annual Report 2022–23*. New Delhi. https://www.meity.gov.in/static/uploads/2024/02/AR_2022-23_English_24-04-23-1.pdf
7. J. M. Patel, B. C. Patel, and M. M. Kayasth, “Advancements in Optical Character Recognition (OCR) for India Scripts: A Review”, *Int. J. Sci. Res. Comp. Sci. Eng.*, vol. 7, no. 1, pp. 23–29, Feb. 2019. <https://ijsrce.isroset.org/index.php/j/article/view/306>

- 466 8. Pal, Umapada, and B. B. Chaudhuri. "Indian script character recognition: a survey." *Pattern*
467 *Recognition* 37, no. 9 (2004): 1887-
468 1899. <https://www.sciencedirect.com/science/article/abs/pii/S003132030400055X>
- 469 9. B.B. Chaudhuri, U. Pal, "A complete printed Bangla OCR system", *Pattern Recognition*
470 31, pp. 531–549,
471 1998. <https://www.sciencedirect.com/science/article/abs/pii/S0031320397000782>
- 472 10. J. Dholakia, A. Yajnik and A. Negi, "Wavelet Feature Based Confusion Character Sets for
473 Gujarati Script," *International Conference on Computational Intelligence and Multimedia*
474 *Applications (ICCIMA 2007)*, Sivakasi, India, 2007, pp. 366-370,
475 <https://ieeexplore.ieee.org/document/4426723?denied=>
- 476 11. Desai, Apurva A. "Gujarati handwritten numeral optical character reorganization through
477 neural network." *Pattern recognition* 43, no. 7, 2582-2589,
478 2010. <https://www.sciencedirect.com/science/article/abs/pii/S0031320310000403>
- 479 12. M. B. Gohil, Dr. A. A. Desai, 2026, A Review on Text Detection and Recognition in
480 Images and Videos Using Machine Learning and Deep Learning Techniques, *International*
481 *Journal Of Engineering Research & Technology (IJERT)* Volume 15, Issue 01, January –
482 2026. [https://www.ijert.org/a-review-on-text-detection-and-recognition-in-images-and-videos-](https://www.ijert.org/a-review-on-text-detection-and-recognition-in-images-and-videos-using-machine-learning-and-deep-learning-techniques-ijertv15is010245)
483 [using-machine-learning-and-deep-learning-techniques-ijertv15is010245](https://www.ijert.org/a-review-on-text-detection-and-recognition-in-images-and-videos-using-machine-learning-and-deep-learning-techniques-ijertv15is010245)
- 484 13. G. Siromoney, R. Chandrasekaran, M. Chandrasekaran, Computer recognition of printed
485 Tamil characters, *Pattern Recognition*, Volume 10, Issue 4, 1978, Pages 243-247, ISSN 0031-
486 3203, <https://www.sciencedirect.com/science/article/abs/pii/0031320378900328>.
- 487 14. Ashwin, T.V., Sastry, P.S. A font and size-independent OCR system for printed Kannada
488 documents using support vector machines. *Sadhana* 27, 35–58 (2002).
489 <https://doi.org/10.1007/BF02703311>
- 490 15. A. Negi, C. Bhagvati and B. Krishna, "An OCR system for Telugu," *Proceedings of Sixth*
491 *International Conference on Document Analysis and Recognition*, Seattle, WA, USA, 2001,
492 pp. 1110-1114, <https://ieeexplore.ieee.org/document/953958?denied=>